# What makes a model?

# Fitting Models

*How do you fit a linear model in R?*
*How many different ways can you think of?*

- `lm` for linear model

- `glm` for generalized linear model (e.g. logistic regression)

- `glmnet` for regularized regression

- `keras` for regression using TensorFlow

- `stan` for Bayesian regression

- `spark` for large data sets

# To specify a model

- Choose a <u>model</u>

- Specify an engine

- Set the mode

## To specify a model

```
1  library(tidymodels)
2  linear_reg()
3  #> Linear Regression Model Specification (regression)
4  #>
5  #> Computational engine: lm
```

# To specify a model

```
1  linear_reg() %>%
2    set_engine("glmnet")
3  #> Linear Regression Model Specification (regression)
4  #>
5  #> Computational engine: glmnet
```

## To specify a model

```
1  linear_reg() %>%
2    set_engine("stan")
3  #> Linear Regression Model Specification (regression)
4  #>
5  #> Computational engine: stan
```
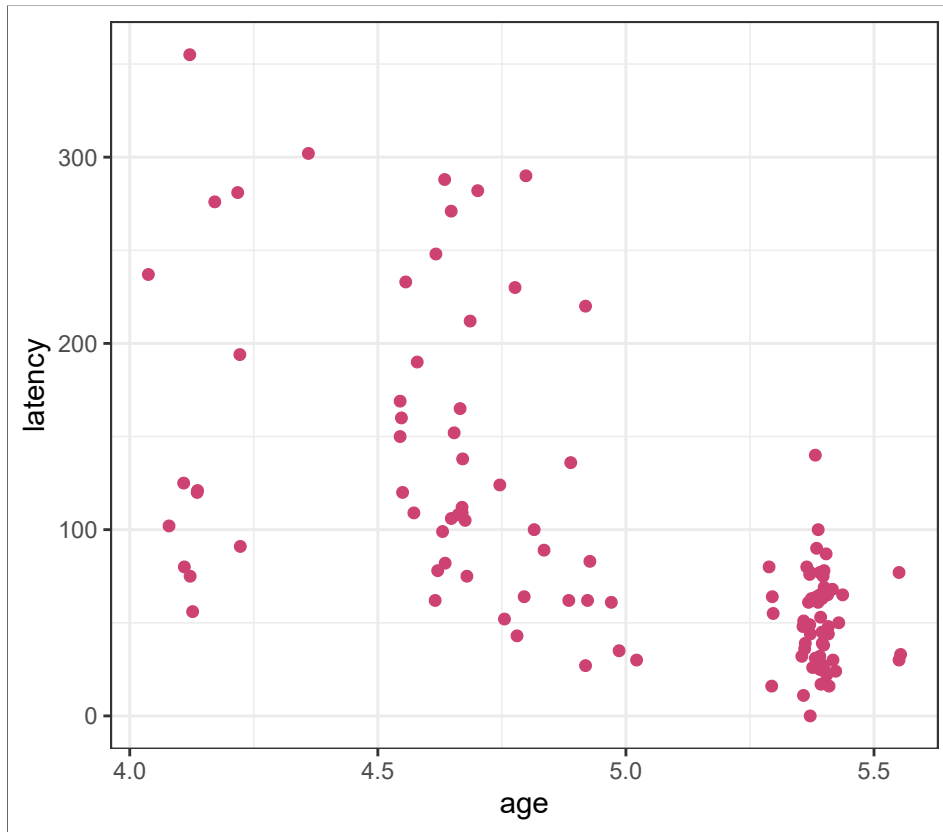
# To specify a model

```
1  decision_tree()
2  #> Decision Tree Model Specification (unknown mode)
3  #>
4  #> Computational engine: rpart
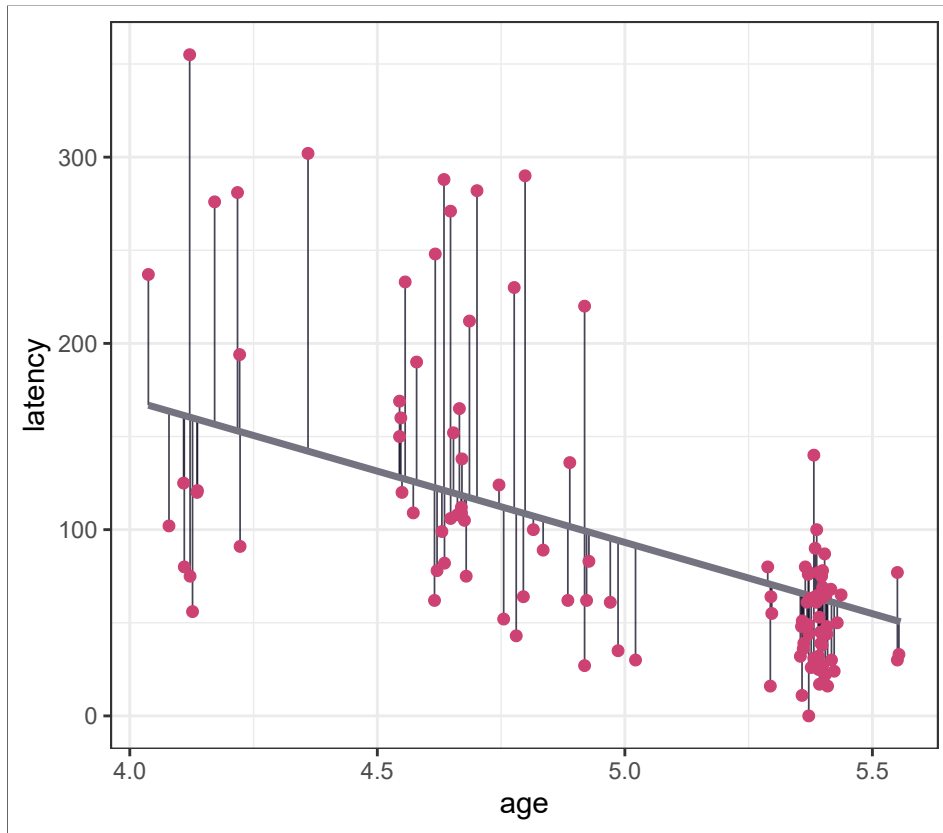```

## To specify a model

```
1  decision_tree() %>%
2    set_mode("regression")
3  #> Decision Tree Model Specification (regression)
4  #>
5  #> Computational engine: rpart
```
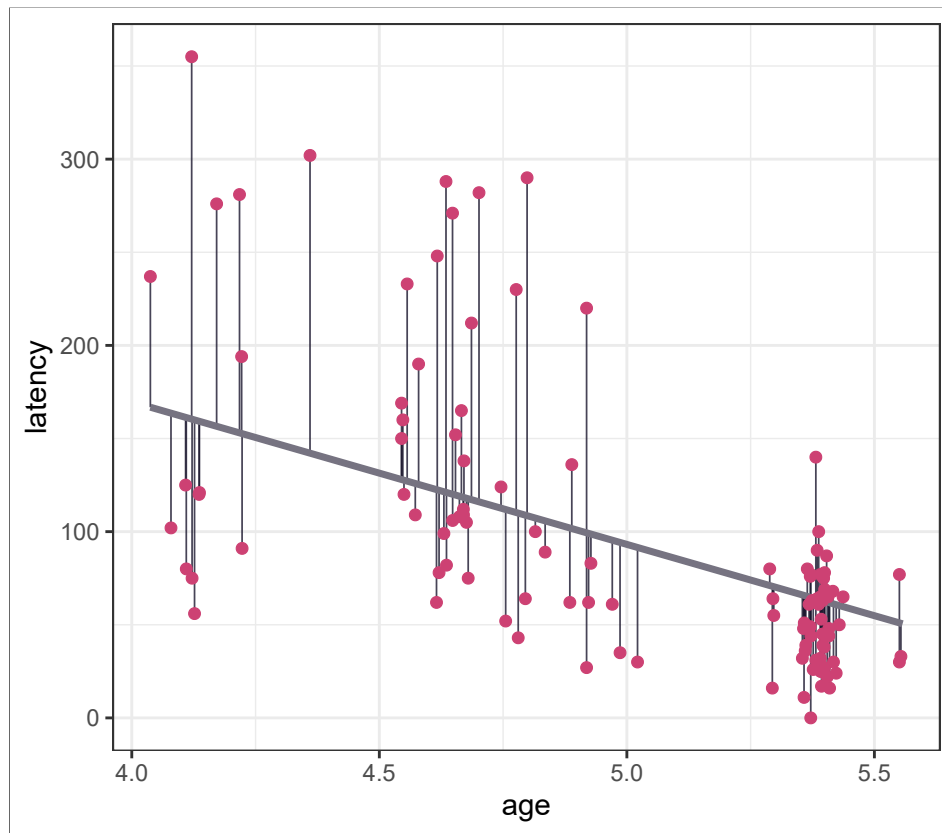
All available models are listed at https://www.tidymodels.org/find/parsnip/

# Linear regression

# Linear regression
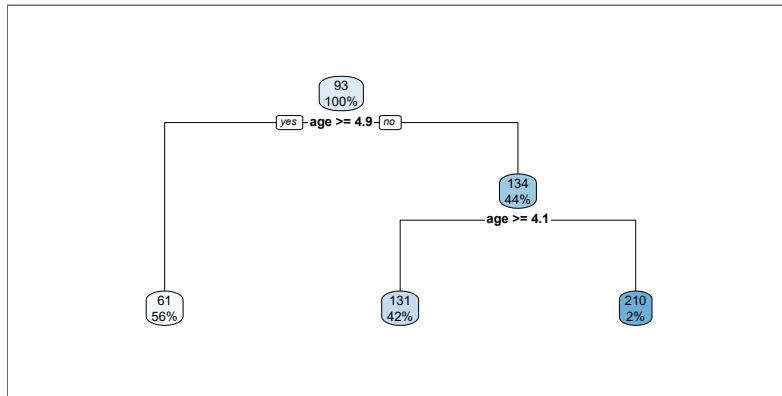
# Linear regression



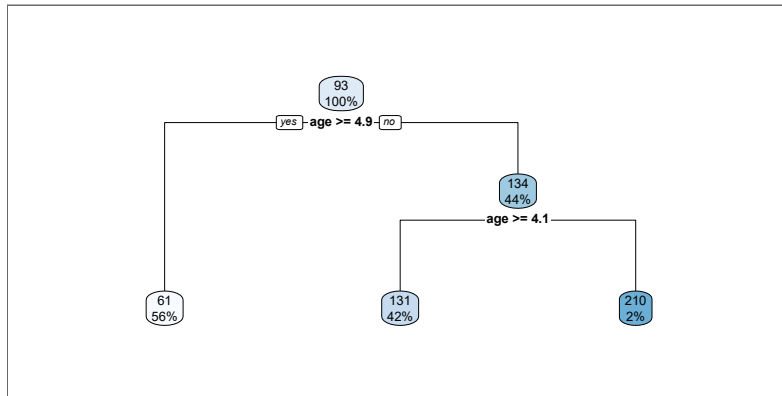- Outcome modeled as linear combination of predictors:

$$\text{latency} = \beta_0 + \beta_1 \cdot \text{age} + \epsilon$$

- Find a line that minimizes the mean squared error (MSE)
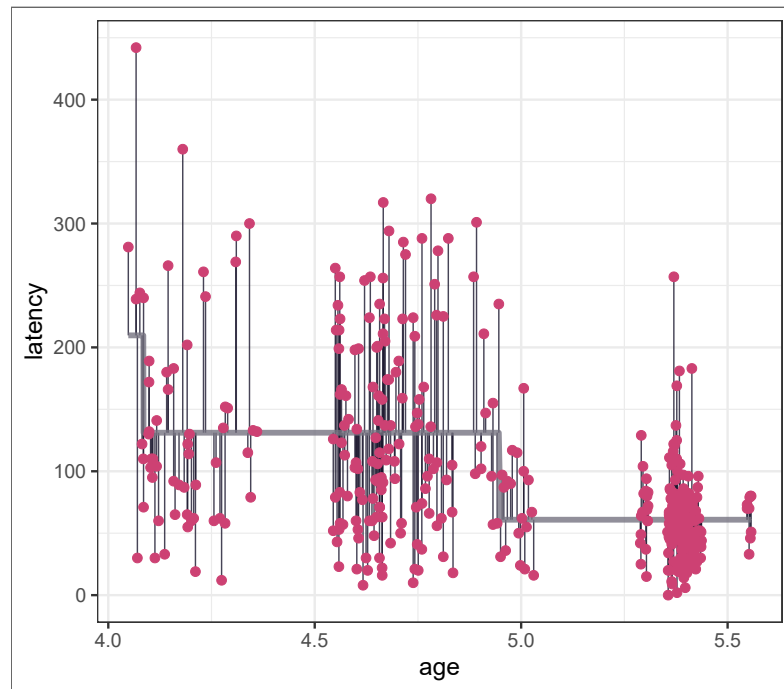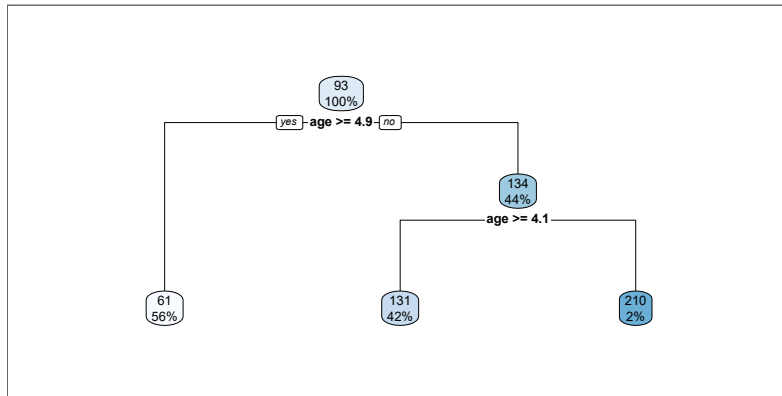
# Decision trees

# Decision trees


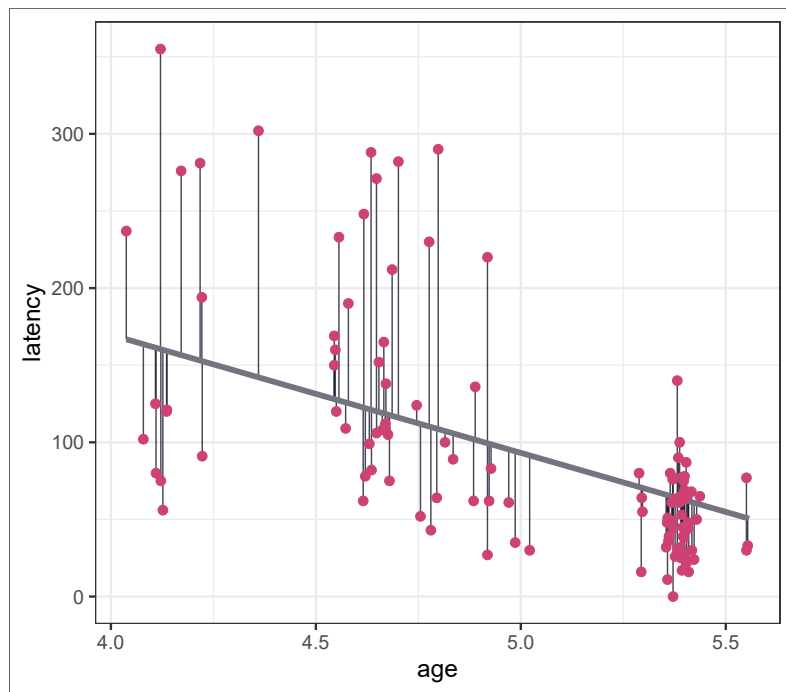
- Series of splits or if/then statements based on predictors

- First the tree *grows* until some condition is met (maximum depth, no more data)

- Then the tree is *pruned* to reduce its complexity
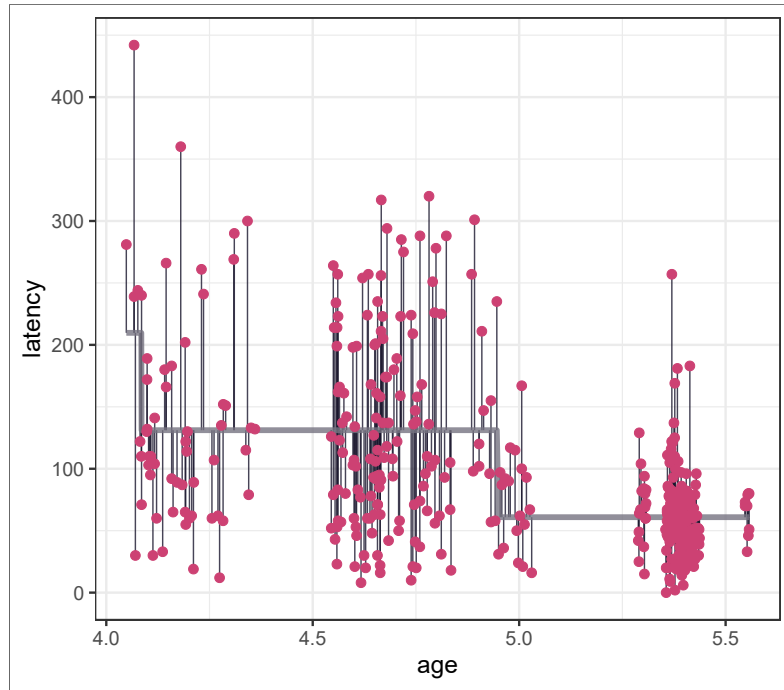
# Decision trees

# All models are wrong, but some are useful!

## Linear regression



## Decision trees

# A model workflow

# Workflows bind preprocessors and models

# What is wrong with this?

# Why a `workflow()`?

- Workflows handle new data better than base R tools in terms of new factor levels

- You can use other preprocessors besides formulas (more on feature engineering tomorrow!)

- They can help organize your work when working with multiple models

- Most importantly, a workflow captures the entire modeling process: `fit()` and `predict()` apply to the preprocessing steps in addition to the actual model fit

# A model workflow

```
 1  tree_spec <-
 2    decision_tree() %>%
 3    set_mode("regression")
 4
 5  tree_spec %>%
 6    fit(latency ~ ., data = frog_train)
 7  #> parsnip model object
 8  #>
 9  #> n= 456
10  #>
11  #> node), split, n, deviance, yval
12  #>       * denotes terminal node
13  #>
14  #>  1) root 456 2197966.00  92.90351
15  #>    2) age>=4.947975 256  252347.40  60.89844
16  #>      4) treatment=control 131   91424.06  48.42748 *
17  #>      5) treatment=gentamicin 125  119197.90  73.96800 *
18  #>    3) age< 4.947975 200 1347741.00 133.87000
19  #>      6) treatment=control 140  986790.70 118.25710
```

# A model workflow

```
 1  tree_spec <-
 2    decision_tree() %>%
 3    set_mode("regression")
 4
 5  workflow() %>%
 6    add_formula(latency ~ .) %>%
 7    add_model(tree_spec) %>%
 8    fit(data = frog_train)
 9  #> ══ Workflow [trained] ═══════════════════════════════════════
10  #> Preprocessor: Formula
11  #> Model: decision_tree()
12  #>
13  #> ── Preprocessor ─────────────────────────────────────────────
14  #> latency ~ .
15  #>
16  #> ── Model ────────────────────────────────────────────────────
17  #> n= 456
18  #>
19  #> node), split, n, deviance, yval
```

# A model workflow

```
1  tree_spec <-
2    decision_tree() %>%
3    set_mode("regression")
4
5  workflow(latency ~ ., tree_spec) %>%
6    fit(data = frog_train)
7  #> ══ Workflow [trained] ══════════════════════════════════════════
8  #> Preprocessor: Formula
9  #> Model: decision_tree()
10 #>
11 #> ── Preprocessor ────────────────────────────────────────────────
12 #> latency ~ .
13 #>
14 #> ── Model ───────────────────────────────────────────────────────
15 #> n= 456
16 #>
17 #> node), split, n, deviance, yval
18 #>       * denotes terminal node
19 #>
```

# Predict with your model

How do you use your new `tree_fit` model?

```r
1  tree_spec <-
2    decision_tree() %>%
3    set_mode("regression")
4
5  tree_fit <-
6    workflow(latency ~ ., tree_spec) %>%
7    fit(data = frog_train)
```
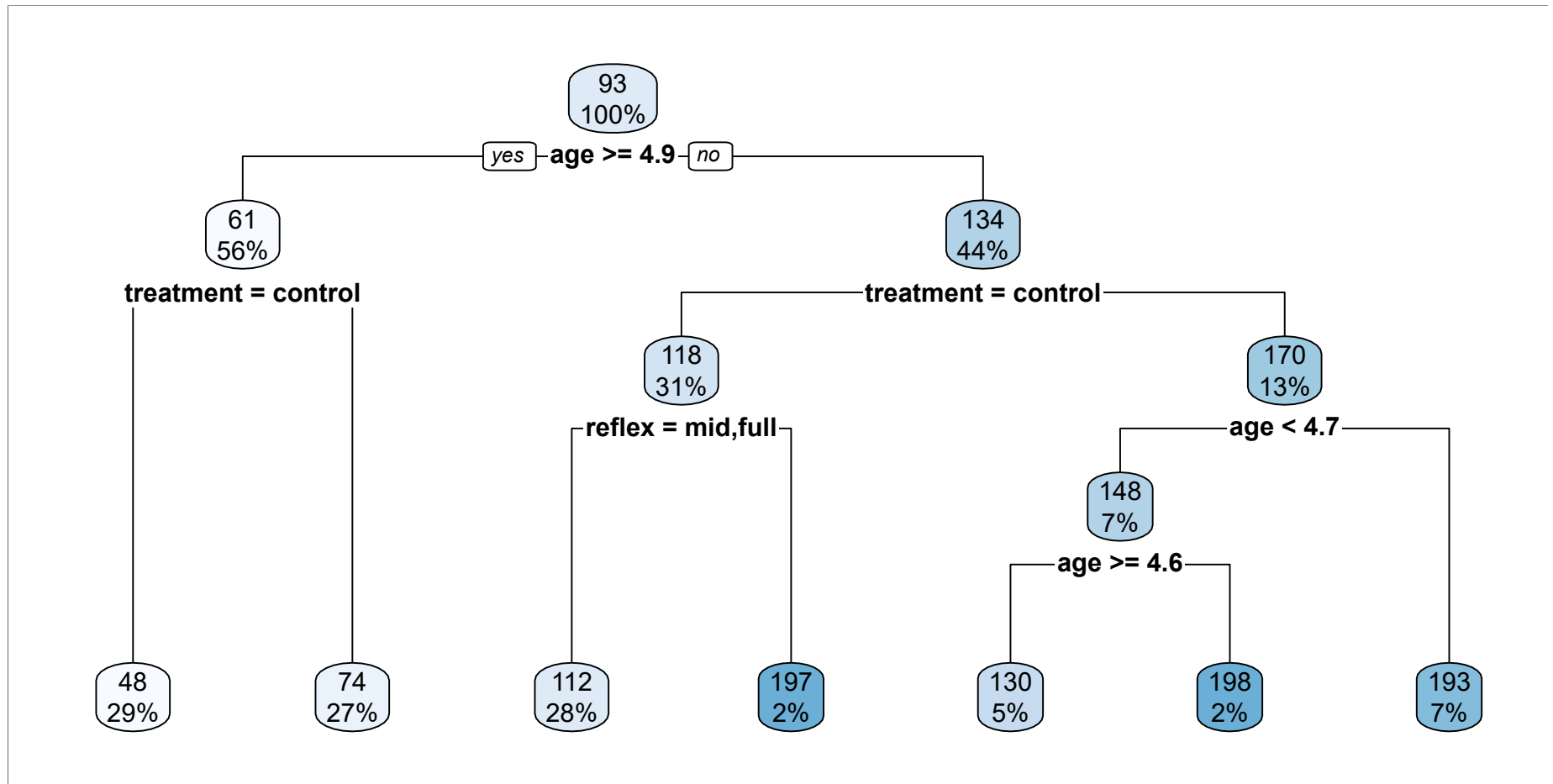
# The tidymodels prediction guarantee!

- The predictions will always be inside a **tibble**

- The column names and types are **unsurprising** and **predictable**

- The number of rows in `new_data` and the output **are the same**

# Understand your model

How do you **understand** your new `tree_fit` model?

# Understand your model

How do you **understand** your new `tree_fit` model?

```
1  library(rpart.plot)
2  tree_fit %>%
3    extract_fit_engine() %>%
4    rpart.plot(roundint = FALSE)
```

You can `extract_*()` several components of your fitted workflow.

# Understand your model

How do you **understand** your new `tree_fit` model?
You can use your fitted workflow for model and/or prediction explanations:

- overall variable importance, such as with the **vip** package

- flexible model explainers, such as with the **DALEXtra** package

Learn more at **https://www.tmwr.org/explain.html**

## Error                                                                          ✕